

The Truth System: Can a System of Lying Processes Stabilize?

Mohamed G. Gouda and Yan Li

Department of Computer Sciences
The University of Texas at Austin
1 University Station (C0500)
Austin, Texas 78712-0233
{gouda, yanli}@cs.utexas.edu

UTCS TR-07-42, August, 2007

Abstract. We introduce a new abstract system, called the truth system. In the truth system, a process deduces a true value, with high probability, from an incoming stream of both true and false values, where the probability that a value in the incoming stream is true is at least 0.6. At each instant, the receiving process maintains at most one candidate of the true value, and eventually the process reaches the conclusion that its candidate value equals, with high probability, the true value. In this paper, we present three versions of the truth system, discuss their properties, and show how to choose their parameters so that their probability of error is small, i.e. about 10^{-6} . The third version, called the stable system, is the most valuable. We employ the stable system as a building block in a stabilizing unidirectional token ring of n processes. When n is small, i.e. about 100 or less, the stable system can be considered error-free and we argue that the resulting token ring is stabilizing with high probability. We simulate the token ring, when n is at most 100, and observe that the ring always stabilizes even though each process lies about its state 40% of the time.

Key words: Distributed Systems, Network Protocols, Self-stabilization

1 Introduction

Faults, that are often assumed to plague the communications between different processes in a distributed system, can be distinguished into natural faults and malicious faults [5]. On one hand, *natural communication faults* are assumed to occur independently of the underlying computation of the distributed system, and are usually assumed to be random. On the other hand, *malicious communication faults* are assumed to occur in the worst possible times for the underlying computation of the distributed system, and are usually assumed to be deliberate and based on complete knowledge (by the adversary) of the underlying computation.

Distributed systems that tolerate natural communication faults are elegant, inexpensive, and practical to implement and use. However, such systems cannot tolerate malicious communication faults if they happen to occur. By contrast, distributed systems that tolerate malicious communication faults are complex, expensive, and sometimes impossible to design.

Well-known examples of natural communication faults and how to tolerate them are as follows.

- *Loss*: Some sent values from one process to another are lost. These faults can be tolerated by making the sending process send the same value repeatedly until the sending process receives an “acknowledgement” from the receiving process [6].
- *Delay*: Some sent values from one process to another are delayed for an unbounded time period before these values are received by their intended receivers. Sometimes, these faults cannot be tolerated as discussed in [4]. This realization has led researchers to adopt weaker model of faults, for example imperfect fault detectors, that can be tolerated [2] and [13].
- *Corruption*: Some sent values are corrupted randomly after they are sent by one process and before they are received by another process. These faults can be detected by adding a checksum to each sent value. In this case, any random corruption of a sent value and its checksum can be detected, with high probability, by the receiving process [11].
- *Topology Change*: The topology of the distributed system changes over time, for example due to the mobility of the processes within the system. Methods for tolerating these “faults” are discussed in [12].
- *Anonymity*: Each sent value does not include the identity of the sending process and is received by an arbitrary process in the distributed system. Methods for tolerating these “faults” are discussed in [1].
- *Modification*: A sent value from one process to another is *modified* before it is received as follows. The value is replaced by any wrong (possibly malicious) value in such a way that the receiving process cannot tell, by examining the received value, that the received value is in fact a false value different from the true value that was sent. For example, if a checksum is attached to the sent value, then a modification fault causes both the value and its checksum to be replaced as follows. The true value is replaced by any false (possibly malicious) value, and the checksum is replaced by the checksum of the false value and so the receiving process cannot tell that the true value and its checksum are replaced by a false value and its checksum. Methods for tolerating modification faults in the context of distributed voting are presented in [10] and [8].

Well-known examples of malicious communication faults and how to tolerate them are discussed in [7] and [9].

In this paper, we present a new method, which we call the truth system, for tolerating modification faults. The truth system is different from the distributed voting systems, discussed in [10] and [8], in several ways. Most notably, the two systems have different objectives. The objective of the truth system is to deduce,

with a predefined high probability, the true value from an *unbounded* stream of both true and false values. In particular, the parameters of the truth system can be set such that the probability of the deduced value being wrong is very small, say 10^{-6} . By contrast, the objective of a distributed voting system is to deduce the value that has the highest probability of being true from a *bounded* stream of both true and false values, but the probability of the deduced value being wrong can be as high as 0.5.

2 Three Versions of the Truth System

The goal of this paper is to introduce a new abstract system, called the truth system, discuss its properties, and show that this system can tolerate, with high probability, modification faults. In the truth system, a process deduces a true value from a mixed stream of both true and false values. The probability, that a value in the mixed stream is true (or false), is at least 0.6 (or at most 0.4, respectively). The process correctly deduces the true value with a high probability, that is around $(1-10^{-6})$.

It is important to explain why we chose the probability, of a value in the mixed stream being true, to be at least 0.6. First, if we chose this probability to be at most 0.5, then no system can deduce the true value with any probability greater than 0. Second, if we chose this probability to be higher than 0.5 but less than 0.6, then as shown at the end of Section 3 the truth system may need up to 1700 values in the input stream to deduce the true value. This is 10 times the number of values needed in the input stream, to deduce the true value, when the probability of a value in the mixed stream being true, is at least 0.6.

The truth system consists of two processes: source and monitor. Process source has an integer state, and process monitor attempts to correctly deduce the integer state of process source.

Periodically, process source sends an integer s to process monitor. With a probability of at least 0.6, the sent s is the true value of the state of process source, and with a probability of less than 0.4, the sent s is an arbitrary integer. Process monitor receives the s integers, one by one, and maintains at most one candidate for the state of process source. Eventually process monitor reaches the conclusion that its maintained candidate for the state of process source equals, with high probability, the true state of process source.

The state of process source can change over time, but we assume that this change occurs at a slow rate. This is because if the state of source is changed at a fast rate, the monitor process may never be able to catch up and deduce the state of source. As shown at the end of Section 3, to ensure that the change rate of the state of process source is slow, we assume that the state of source is not changed until process source has sent out this state 170 or more times.

Our presentation of the truth system in this paper consists of three steps. In each step, we present a version of the truth system and discuss its properties. We then point out some problem with this version and so clear the way for the

next version that is to be presented in the next step, and so on. The version presented in the third (and last) step has no problems, as far as we can tell.

In the first step, we present a version of the truth system where the monitor process terminates as soon as it concludes that its maintained candidate for the state of process source equals, with high probability, the true state of process source. The problem of this truth system is that the monitor process deduces exactly one true state of process source, even if the state of process source is changed many times afterwards. We refer to this truth system as the *one-shot system*.

In the second step, we modify the one-shot system to make the monitor process continue to operate indefinitely, even after it concludes that its maintained candidate for the state of process source equals, with high probability, the true state of process source. We refer to this truth system as the *continuing system*. The problem with the continuing system is that the conclusion reached by process monitor (that its candidate for the state of process source equals, with high probability, the true state of source) is not stable, but it can fluctuate wildly over time, even when the true state of process source remains fixed for a long time period.

In the third step, we modify the continuing system to ensure that, when the true state of process source remains fixed for a long time period, the conclusion reached by process monitor (that its candidate for the state of process source equals, with high probability, the true state of source) remains stable over time. We refer to this truth system as the *stable system*.

3 The One-Shot System

In this section, we present our first version of the truth system, called the one-shot system. In the one-shot system, as soon as process monitor concludes, that its maintained candidate for the state of process source equals, with high probability, the true state of process source, process monitor terminates.

Process source in the one-shot system is specified in Protocol 1. This process has one action that it executes over and over, since the guard of the action is **true**. During each execution of the action, process source sends an integer s to process monitor. With probability $p/100$, the sent s is the value of input state, and with probability $(100-p)/100$, the sent s can be any integer, where p is an input of process source.

Input p can be regarded as the *probability of process source telling the truth*, and input *state* can be regarded as the *true state of process source*. From the received s integers, process monitor is expected to deduce the true value of input *state* in process source. It is straightforward to show that if the value of p is in the range 0..50, then process monitor can never deduce the true value of *state*. At the end of Section 3, we argue that if the value of p is in the range 51..59, then process monitor may need to receive up to 1700 s integers (instead of 170) in order to deduce the true value of *state* with high probability. That is why we

Protocol 1 process source

```

const   pmin : 60..99
input   p    : pmin..99
          state : integer {state of source}
variable r    : 0..99    {random number}
          s    : integer {sent state}
begin
  true →
    r := random
    if r ≥ p then
      s := any           {assign malicious value}
    else
      s := state
    end
    send s to monitor
end

```

specified the value of p to be in the range $pmin..99$, where $pmin$ is a constant whose value is in the range $60..99$.

Both p and $state$ are inputs to process source. Thus, their values can be changed over time by an outside agent. As shown at the end of Section 3, we assume that once the value of state is changed, the value of state remains fixed until process source executes its action 170 or more times.

Process monitor in the one-shot system is specified in Protocol 2. Process monitor has three variables: c , st and s . Variable c is a counter whose value is in the range $0..cmax$, where $cmax$ is a constant of process monitor. Variable st stores the latest candidate for the state of process source. Variable s stores the latest received integer from process source. The value of counter c indicates whether process monitor can conclude that the current value of st equals, with high probability, the value of state in process source. Process monitor reaches this conclusion when, and only when, the value of counter c is $cmax$.

Process monitor has only one action that is executed each time the process receives an integer s from process source. When an integer s is received, process monitor checks the value of its counter c . If $c = 0$, then variable st is assigned s and counter c is assigned 1. If $c > 0$ and st is different from the received s , then c is decreased by 1. If $c > 0$ and st equals the received s , then c is increased by 1 (provided that c does not exceed its maximum value $cmax$). Then process monitor compares the values of c with $cmax$. If $c = cmax$, then process monitor concludes that the current value of its variable st equals, with high probability, the value of $state$ in process source.

To complete the specification of process monitor, we need now to compute the value of constant $cmax$ in monitor. The value of $cmax$ should be chosen such that the probability of error of the one-shot system is kept small, say around 10^{-6} .

Protocol 2 process monitor of the one-shot system

```

const   pmin : 60..99  {same as pmin in source}
          cmax : integer
variable c   : 0..cmax {counter, init. 0}
          s   : integer {received state}
          cs  : integer {candidate state}
begin
  rcv s from source  $\rightarrow$ 
    if c = 0 then
      c := 1
      cs := s
    else if cs  $\neq$  s then
      c := c - 1
    else
      c := min(c+1, cmax)
    end
    if c = cmax then
      {conclude: cs = state} terminate
    end
end

```

The *probability of error*, denoted $p(\text{error})$, of the one-shot system is the probability that starting from its initial global state where $c = 0$, the system reaches a global state where $c = cmax$ and $cs \neq state$.

In Theorem 1 below, we give a formula that describes the relationship between $pmin$, $cmax$, and $p(\text{error})$ for the one-shot system. Our proof of this theorem is based on two simplifying assumptions. First, we assume that the state of process source does not change over time. This assumption is acceptable given our understanding that the change rate of the state of process source is slow anyway. Second, we assume that whenever process source sends an arbitrary integer s to process monitor, process source always sends the same integer that is different from the state of process source. This assumption represents the worst case scenario that assigns $p(\text{error})$ its highest value. We adopt these two assumptions in proving all the theorems in this paper. All the proofs are shown in the Appendix at the end of the paper.

Theorem 1 (pmin, cmax, and p(error) for the one-shot system).

$$p(\text{error}) = \frac{(1 - pmin)^{cmax}}{(1 - pmin)^{cmax} + (pmin)^{cmax}} . \quad \square$$

For many applications, it is reasonable to expect that $p(\text{error})$ should be around 10^{-6} . In this case, we can use the formula in Theorem 1 to produce the relationship between $pmin$ and $cmax$, for the one-shot system, shown in Table 1.

An execution *step* of the one-shot system consists of two parts. First, process source executes its (sending) action, then process monitor executes its (receiving) action.

Table 1.

$pmin$	$cmax$
0.6	34
0.7	16
0.8	10
0.9	6

If $pmin$ in this system is 0.6 and the value of variable cs in process monitor is the correct state of process source, then in each step of the system, counter c in process monitor is incremented by 1 with probability 0.6, and is decremented by 1 with probability 0.4. In other words, each step of the system increments counter c by 0.2 on the average. Thus, the system needs to execute $cmax/0.2$ steps on the average before counter c reaches its maximum value $cmax$ and the system terminates. Because $cmax$ in this system is 34 from Table 1, the system needs to execute 170 steps before it terminates.

If we choose $pmin$ in this system to be 0.51, and assume that $cmax$ remains 34 (rather than being increased in value as it should), and follow the same analysis in the previous paragraph, we conclude that the system will execute on average $(34/0.02)=1700$ steps before it terminates. In other words, choosing $pmin$ to be in the range from 0.51 to 0.59 can lead to (sometimes substantial) increase in the number of steps to be executed. This should explain our choice of $pmin$ to be at least 0.6.

The above analysis for computing the average number of steps that need to be executed by the one-shot system before it terminates is based on the assumption that the state of process source does not change during execution. This is an important assumption; for instance, if the state is changed at least once every 5 execution steps, the one-shot system may never terminate. This should explain our above requirement that the state of process source remains fixed for the duration of 170 steps.

4 The Continuing System

The problem of the one-shot system is that process monitor terminates as soon as it concludes that the value of its cs variable equals, with high probability, the value of input state in process source. Thus, monitor cannot observe any change in the state of process source. To remedy this problem, we modify process monitor such that the process continues to execute indefinitely. This modification is achieved by replacing statement **terminate** by a statement **skip** in the action of process monitor. We refer to the resulting system as the continuing system.

Because the continuing system is nonterminating, the initial state of the system, where $c = 0$, is irrelevant. Rather, we define the *probability of error* $p(error)$ of the continuing system as the steady state probability that the system is in a global state where $c = cmax$ and $cs \neq state$. The following theorem gives a formula that describes the relationship between $pmin$, $cmax$, and $p(error)$ for the continuing system.

Theorem 2 (pmin, cmax, and p(error) for the continuing system).

$$p(error) = \frac{(1 - pmin)^{2 \times cmax}}{(1 - pmin)^{2 \times cmax} + (pmin)^{2 \times cmax}} . \quad \square$$

Assuming that $p(error)$ is around 10^{-6} , we can use the formula in Theorem 2 to produce the relationship between $pmin$ and $cmax$, for the continuing system, shown in Table 2. Notice that the $cmax$ values in the one-shot system, shown in Table 1, are twice the $cmax$ values in the continuing system, shown in Table 2.

The continuing system has an interesting problem. Even if the state of process source remains fixed for a long time period T , the value of counter c in process monitor can fluctuate during period T between $c < cmax$ (when process monitor cannot conclude that $cs = state$) and $c = cmax$ (when process monitor can conclude that $cs = state$). This observation suggests the following definition.

The *probability of no conclusion*, denoted $p(no-conclusion)$, of the continuous system is the steady state probability that the system is in a global state where $c < cmax$. The following theorem gives a formula for computing $p(no-conclusion)$ as a function of $pmin$ and $cmax$.

Theorem 3 (p(no-conclusion) for the continuing system).

$$p(no-conclusion) = \frac{\sum_{i=1}^{2 \times cmax - 1} \left(\frac{1 - pmin}{pmin}\right)^i}{\sum_{j=0}^{2 \times cmax} \left(\frac{1 - pmin}{pmin}\right)^j} . \quad \square$$

Using the formula in Theorem 3 and the values of $pmin$ and $cmax$ in Table 2, we can compute the values of $p(no-conclusion)$, for the continuing system, as shown in Table 2.

Table 2.

$pmin$	$cmax$	$p(no-conclusion)$
0.6	17	0.667
0.7	8	0.429
0.8	5	0.250
0.9	3	0.111

From the first row in Table 2, if $pmin$ and $cmax$ for the continuing system are 0.6 and 17 respectively, then $p(no-conclusion)$ for this system is 0.667. This means that even if the state of process source remains fixed for a long time period T , process monitor cannot conclude that $cs = state$ for 66.7% of the time during period T . Clearly, this is not acceptable and a modification of the continuing system is in order. In the next section, we describe how to modify the continuing system to remedy this problem. We refer to the modified system as the stable system.

5 The Stable System

The stable system is obtained from the continuing system, discussed in the previous section, by making the following two modifications to process monitor (in the continuing system). First, a new variable named *ss* is added to process monitor. Variable *ss* stores the latest stable estimate (by process monitor) of the state of process source. Second, the last **if**-statement in the action of process monitor is modified to become as shown in Figure 1. (The first **if**-statement in the action of process monitor remains unchanged.)

```

if  $c = cmax$  then
    {conclude:  $cs = state$ }  $ss := cs$ 
end
    
```

Fig. 1.

The *probability of error* $p(error)$ of the stable system is defined as the steady state probability that the system is in a global state where $ss \neq state$. The following theorem describes the relationship between $pmin$, $cmax$, and $p(error)$ for the stable system.

Theorem 4 (pmin, cmax, and p(error) for the stable system). *Given that $p(error)$ for the stable system is around 10^{-6} , we have the following relationship between $pmin$ and $cmax$ for the stable system.*

Table 3.

$pmin$	$cmax$
0.6	20
0.7	9
0.8	5
0.9	4

□

A *step* of the stable system consists of two parts. First, process source executes its action once. Second, process monitor executes its action once. The *convergence span* of the stable system is the average number of steps that needs to be executed by the stable system in order to change the global state of the system from one where $c = cmax$ and $ss \neq state$ to one where $c = cmax$ and $ss = state$. The following theorem gives an approximate formula for computing the convergence span of the stable system.

Theorem 5 (convergence span of the stable system).

$$\text{convergence span} \approx \frac{2 \times cmax}{2 \times pmin - 1} . \quad \square$$

Using the formula in this theorem and the values of $pmin$ and $cmax$ from Table 3, we compute the convergence span of the stable system as shown in Table 4.

$pmin$	$cmax$	convergence span
0.6	20	200
0.7	9	45
0.8	5	17
0.9	4	10

Table 4.

6 A Stabilizing Token Ring

In this section, we discuss how to employ the stable system, presented in the previous section, as a building block in constructing a stabilizing unidirectional token ring of up to 100 processes, where each process can lie about its state at most 40% of the time. We start our discussion by presenting a unidirectional token ring, in Protocol 3, where processes do not lie about their states. Note that this ring is similar to Dijkstra’s token ring in [3] with two exceptions. First, the execution of this ring is synchronous, whereas the execution of Dijkstra’s token ring is asynchronous. Second, this ring uses message passing primitives whereas Dijkstra’s token ring uses shared memory primitives.

In this ring, each process $p[i]$ has two variables, s and ss , where variable s stores the latest state that $p[i]$ has sent or received, and variable ss stores the state of $p[i]$. Each $p[i]$ also has two actions: a sending action where $p[i]$ sends its own state to $p[i+1 \bmod n]$, and a receiving action where $p[i]$ receives the state of $p[i-1 \bmod n]$ then modifies its own state based on the received state.

A *global state* of this ring is defined by a value for each ss variable in the ring. (This means that the s variables are not considered part of the global state of the ring.)

A *transition* of this ring is a pair (S, S') of global states of the ring such that if the ring is in a global state S and a “step” is executed, then the ring becomes in a global state S' . Executing a step in the ring consists of two parts. First, each process in the ring executes its sending action, then each process in the ring executes its receiving action. Thus, each process in the ring ends up executing both its (sending and receiving) actions in a step.

A *computation* of this ring is an infinite sequence $S.0, S.1, \dots$ of global states of the ring such that each pair $(S.i, S.(i+1))$ of consecutive states in the sequence is a transition of the ring.

Protocol 3 process $p[i : 0..n-1]$ in the original token ring

```

variable s : 0..n-1 {sent/received state}
           ss : 0..n-1 {state}
begin
  true →
    s := ss
    send s to p[i+1 mod n]

||   rcv s from p[i-1 mod n] →
      if i > 0 then
        ss := s
      else if ss = s then
        ss := ss + 1 mod n
      end

end

```

It is straightforward to show that each computation of this ring reaches a legitimate global state where the values of all the ss variables are equal after at most $2n$ transitions, and so this ring is stabilizing.

Clearly, stabilization of the ring in Protocol 3 depends heavily on the fact that the ring processes do not lie when they send their states to other processes. To allow the ring processes to lie about their states, 40% of the time, and still retain the stabilization of the ring, we employ the stable system, discussed in the previous section, in constructing the new ring. Specifically, each process $p[i]$ in the ring is modified to act as a source when $p[i]$ sends a state s to process $p[i+1 \bmod n]$, and act as a monitor when $p[i]$ receives a state s from process $p[i-1 \bmod n]$. The new ring is shown in Protocol 4.

A *global state* of the new ring is defined by a value for each ss variable in the ring. (This means that none of the other variables, namely r , c , s , and cs , is considered part of the global state of the new ring.)

A *transition* of the new ring is a pair (S, S') of the global states of the new ring such that if the ring is in a global state S and a “step” is executed, then the ring becomes in a global state S' . Executing a step in the new ring consists of two parts. First, each process in the ring executes its sending (or source) action. Second, each process in the ring executes its receiving (or monitor) action.

A *computation* of the new ring is an infinite sequence $S.0, S.1, \dots$ of global states of the ring such that each pair $(S.i, S.(i+1))$ of consecutive states in the sequence is a transition of the new ring.

The new ring can be viewed as consisting of n stable systems, and each ss variable in the ring can be viewed as belonging to the monitor of one of those stable systems. For the new ring to stabilize, if it is to stabilize, each ss variable needs to be assigned a value around $2 \times n$ times. From Table 4 and given that p_{\min} is 0.6, a stable system needs to execute on average 200 steps in order to assign a value to its ss variable. Therefore, for the new ring to stabilize, if it is to

Protocol 4 process $p[i : 0..n-1]$ in the new token ring

```

const   pmin : 60..99 {pmin = 60}
          cmax : integer {cmax = 20}
variable r   : 0..99 {random number}
          c   : 0..cmax {counter, init. 0}
          s   : 0..n-1 {sent/received state}
          cs  : 0..n-1 {candidate state}
          ss  : 0..n-1 {stable state}

begin
  true →
    r := random
    if r ≥ pmin then
      s := any
    else
      s := ss
    end
    send s to p[i+1 mod n]

  || rcv s from p[i-1 mod n] →
    if c = 0 then
      c := 1
      cs := s
    else if cs ≠ s then
      c := c - 1
    else
      c := min(c+1, cmax)
    end
    if c = cmax then
      if i > 0 then
        ss := cs
      else if ss = cs then
        ss := ss + 1 mod n
      end
    end
end

```

stabilize, each of the n stable systems in the ring needs to execute around $n \times 400$ steps. By choosing n to be relatively small, say 100, each stable system in the ring needs to execute a small number of steps, around 40000 steps in order for the ring to stabilize. Because the probability of error of a stable system is very small, around 10^{-6} , it is reasonable to assume that whenever any ss variable is assigned a value, in the first 40000 transitions of a computation, it is assigned a correct value. We refer to this assumption as the *no-use-lying* assumption.

Now consider a computation $S.0, \dots, S.40000, \dots$ of the new ring. Under the no-use-lying assumption, whenever an ss variable is assigned a value in the first 40000 transitions of this computation, it is assigned a correct value. Therefore,

the global state S.40000 in this computation is a legitimate state, with high probability. Therefore, the new ring is stabilizing, with high probability.

The probabilistic stabilization of the new ring depends heavily on the validity of the no-use-lying assumption. To check the validity of this assumption, we have run 100 simulations of the new ring. For each simulation, we chose n (the number of processes in the ring) to be 100, the initial global state of the ring to be random, and the wrong state that each process sends in place of its correct state to be 0. We observed that each simulation has stabilized to a legitimate global state, after no more than 32440 transitions. These simulation results justify our adoption of the no-use-lying assumption.

7 Concluding Remarks

The truth system is a building block that can be employed in a distributed system to ensure that the system performs its intended function, with high probability, even if up to 40% of the sent values by each process in the system are completely arbitrary. In this paper, we presented three versions of the truth system: the one-shot system, the continuing system, and the stable system. We also compared the properties of these three versions and concluded that the stable system is superior to the other two. Finally we showed how to employ the stable system in a unidirectional token ring so that the ring performs its intended function even if up to 40% of the values sent by each process in the ring are arbitrary.

This paper suggests a number of interesting problems that merit further research. First, are there interesting versions of the truth system other than those discussed in this paper? Second, are there algorithms that take a distributed system that performs a function f under the assumption of perfect communication and produce a distributed system that employs a version of the truth system as a building block and performs function f , with high probability, under the assumption that up to 40% of the values sent by each process in the system are arbitrary. Third, are there effective methods to compute the probability of error and the convergence span for a distributed system where a version of the truth system is employed as a building block?

Acknowledgments. The work of M. G. Gouda is supported in part by a grant 0520250 from the National Science Foundation.

References

1. D. Angluin, J. Aspnes, D. Eisenstat, and E. Ruppert. On the power of anonymous one-way communication. In *Principles of Distributed Systems; 9th International Conference, OPODIS*, Pisa, Italy, December 2005.
2. T. D. Chandra and S. Toueg. Unreliable failure detectors for asynchronous systems. *Journal of the ACM*, 43(2), March 1996.
3. E. W. Dijkstra. Self-stabilization systems in spite of distributed control. *CACM*, November 1974.

4. M. J. Fischer, N. A. Lynch, and M. S. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 32(2), April 1984.
5. M. J. Fischer and M. Merritt. Appraising two decades of distributed computing theory research. *Distributed Computing*, 16(2-3), 2003.
6. A. Herzberg and S. Kutten. Early detection of message forwarding faults. *SIAM Journal of Computing*, August-October issue, 2000.
7. E. N. Hoch, D. Dolev, and A. Daliot. Self-stabilizing byzantine digital clock synchronization. In *8th International Symposium on Stabilization, Safety, and Security of Distributed Systems*, Lecture Notes in Computer Science 4280. Springer, 2006.
8. A. Kumar and K. Malik. Voting mechanisms in distributed systems. *IEEE Transactions on Reliability*, 40(5), December 1991.
9. M. R. Malekpour. Byzantine-fault tolerant self-stabilizing protocol for distributed clock synchronization systems. In *8th International Symposium on Stabilization, Safety, and Security of Distributed Systems*, Lecture Notes in Computer Science 4280. Springer, 2006.
10. M. Paquette and A. Pelc. Optimal decision strategies in byzantine environments. *Journal of Parallel and Distributed Computing*, 66(3), March 2006.
11. W. W. Peterson and D.T. Brown. Cyclic codes for error detection. In *Proceedings of the IRE*, volume 49, 1961.
12. J. E. Walter, J. L. Welch, and N. H. Vaidya. A mutual exclusion algorithm for ad hoc mobile networks. *Wireless Networks*, 2001.
13. J. Yang, E. Gafni, and G. Neiger. Structured derivations of consensus algorithms for failure detectors. In *PODC*, Puerto Vallarta, Mexico, 1998.

8 Appendix: Proofs of Theorems

8.1 Proof of Theorem 1

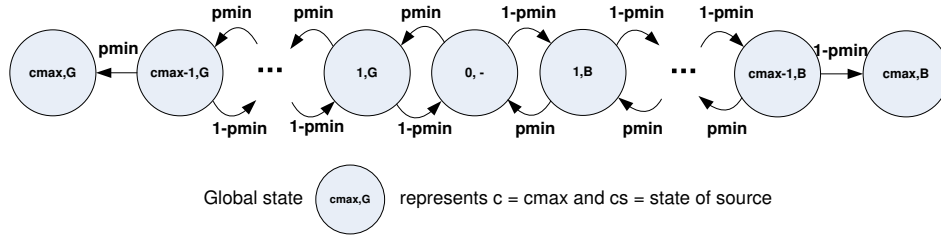


Fig. 2. State transition graph for the one-shot system

The state transition graph for the one-shot system is shown in figure 2. There are two elements in each global state in the state transition graph. The first element is the current counter value c , and the second element represents whether the candidate state cs maintained in process monitor equals $state$ of process source. If $cs = state$ then we say cs is good(G), otherwise, cs is bad(B). When $c = 0$, cs can be either good(G) or bad(B).

The probability of error when counter c reaches $cmax$ is the probability from the start state $(0, -)$ to the end state $(cmax, B)$. Let p_G denotes the probability from start state $(0, -)$ to end state $(cmax, G)$ and p_B denotes the probability from start state $(0, -)$ to end state $(cmax, B)$. Thus we have $p_B = p(error)$. Since process monitor will terminate when $c = cmax$, it can only terminate at one of the two end states, which is either $(cmax, G)$ or $(cmax, B)$. So the following formula holds:

$$p_G + p_B = 1.0 . \tag{1}$$

p_B is the sum of probability of each path from start state $(0, -)$ to end state $(cmax, B)$. Because of the symmetric property of the state transition graph, each global state with cs being bad(B) has a one-to-one corresponding state with cs being good(G). Thus, for every path from $(0, -)$ to $(cmax, B)$, called *original path*, we can construct exactly one corresponding path from $(0, -)$ to $(cmax, G)$, which is called a *mirrored path*. Furthermore, for each loop in the mirrored path, the probability of that loop is the same as the one in the original path. We call this property *loop symmetry*. For each path from $(0, -)$ to $(cmax, B)$, no matter how many loops it has, the probability must have the multiplication factor $(1 - pmin)^{cmax}$ since the path finally need to go out of the loops and reach final state $(cmax, B)$. Similarly, for each path from $(0, -)$ to $(cmax, G)$, the probability must have the multiplication factor $pmin^{cmax}$. And the loops in the mirrored paths have the same probability as that in the original paths

because of *loop symmetry*. Thus, we have the following formula:

$$\frac{p_B}{p_G} = \frac{(1 - pmin)^{cmax}}{(pmin)^{cmax}} . \quad (2)$$

From equation 1 and 2, we can get

$$p_B = \frac{(1 - pmin)^{cmax}}{(1 - pmin)^{cmax} + (pmin)^{cmax}} . \quad (3)$$

Since $p(error) = p_B$, the theorem holds.

8.2 Proof of Theorem 2

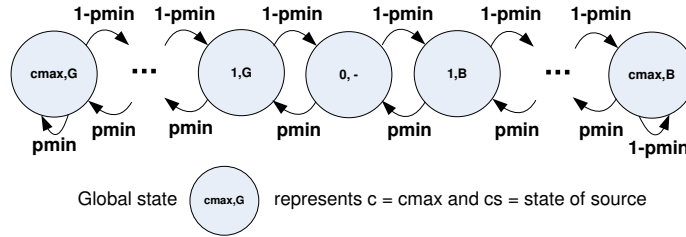


Fig. 3. State transition graph for the continuing system

Since the starting state has no effect in steady state, process monitor of the continuing system can be described as a birth-death process as shown in Figure 3. The probability of giving a birth at any time t is $1 - pmin$ while the probability of death at any time t is $pmin$. Let $i, i = 0, \dots, 2 \times cmax$ denote the i th global state from left to right in Figure 3. And let $\pi_i, i = 0, \dots, 2 \times cmax$ denotes the probability at global state i . The stationary distribution of this birth-death process is:

$$\pi_{i+1} = \frac{1 - pmin}{pmin} \times \pi_i, \text{ for } i = 0, \dots, 2 \times cmax . \quad (4)$$

$$\sum_{i=0}^{2 \times cmax} \pi_i = 1 . \quad (5)$$

From equation 4, we can get:

$$\pi_i = \left(\frac{1 - pmin}{pmin} \right)^i \times \pi_0, \text{ for } i = 0, \dots, 2 \times cmax . \quad (6)$$

By the definition of $p(error)$ in the continuing system, we have

$$p(error) = \frac{\pi_{2 \times cmax}}{\pi_0 + \pi_{2 \times cmax}} . \quad (7)$$

Substituting 6 in 7, we can get:

$$p(error) = \frac{(1 - pmin)^{2 \times cmax}}{(1 - pmin)^{2 \times cmax} + (pmin)^{2 \times cmax}} . \quad (8)$$

Thus the theorem holds.

8.3 Proof of Theorem 3

By the definition of $p(no-conclusion)$ in the continuing system, we have

$$p(no-conclusion) = \sum_{i=1}^{2 \times cmax - 1} \pi_i . \quad (9)$$

From equation 4 and 5, we have

$$\pi_0 = \frac{1}{\sum_{i=0}^{2 \times cmax} \left(\frac{1-pmin}{pmin}\right)^i} . \quad (10)$$

Substituting 6 and 10 in 9, we can get:

$$p(no-conclusion) = \frac{\sum_{i=1}^{2 \times cmax - 1} \left(\frac{1-pmin}{pmin}\right)^i}{\sum_{j=0}^{2 \times cmax} \left(\frac{1-pmin}{pmin}\right)^j} . \quad (11)$$

8.4 Proof of Theorem 4

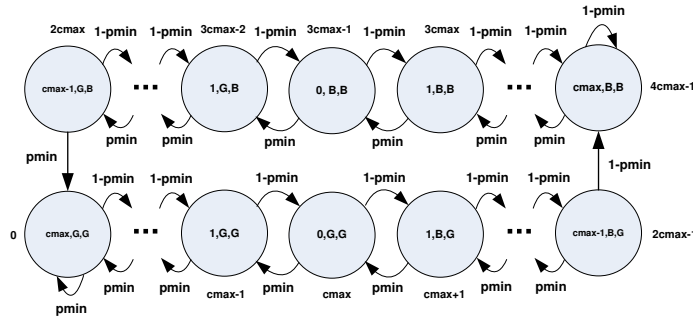


Fig. 4. State transition graph for the stable system

The state transition graph for the stable system is shown in Figure 4, which basically is a Markov chain. Each global state includes three elements as follows. The first element is the current counter value c of process monitor. The second element represents whether the candidate value cs maintained in process monitor equals the *state* of process source. We use G (*Good*) to represent $cs = state$ and use B (*Bad*) to represent $cs \neq state$. The third element represents whether the stable state ss maintained in process monitor equals the *state* of process source. Similar to the second element, we use G (*Good*) to indicate that $ss = state$ and use B (*Bad*) otherwise. The initial state is $(0, B, B)$ which means in process monitor, counter c is zero and neither cs nor ss equals the *state* of process source. As stated in Figure 1, ss can only be changed to cs when the counter value c is equal to $cmax$, and ss will remain unchanged even when the candidate value cs changes but counter value c is less than $cmax$.

Let $i, i = 0, \dots, 4 \times cmax - 1$ denote the i th global state from left to right as shown in Figure 4. And let $\pi_i, i = 0, \dots, 4 \times cmax - 1$ denotes the probability at the i th global state in equilibrium and let $\pi = [\pi_0, \pi_1, \dots, \pi_{4 \times cmax - 1}]$. By the definition of $p(error)$ in the stable system, we have the following equation:

$$PE = \sum_{i=2 \times cmax}^{4 \times cmax - 1} \pi_i . \quad (12)$$

Since the state transition graph forms a Markov chain, we have the following equations:

$$\sum_{i=0}^{4 \times cmax - 1} \pi_i = 1.0 . \quad (13)$$

$$\pi \times P = \pi . \quad (14)$$

where P is the transition probability matrix of the Markov chain.

$$P(i, j) = \begin{cases} 1 - tr & \text{if } \alpha \\ tr & \text{if } \beta \\ 0 & \text{otherwise} \end{cases}$$

The boolean conditions α and β are defined as follows:

$$\begin{aligned} \alpha &= (0 \leq i < 4 \times cmax - 2 \wedge i \neq 2 \times cmax - 1 \wedge j = i + 1) \\ &\quad \vee (i = 2 \times cmax - 1 \wedge j = 4 \times cmax - 1) \\ \beta &= (1 \leq i < 4 \times cmax - 1 \wedge i \neq 2 \times cmax \wedge j = i - 1) \\ &\quad \vee (i = 2 \times cmax \wedge j = 0) \end{aligned}$$

By solving equations 13 and 14, we can get the probability at the i th global state π_i for $i = 0, \dots, 4 \times cmax - 1$. Thus we can compute $p(error)$ using equation 12. When $pmin = 0.6$ and $cmax=20$, $p(error) = 1.2 \times 10^{-6}$; when $pmin = 0.7$ and $cmax=9$, $p(error) = 2.3 \times 10^{-6}$; when $pmin = 0.8$ and $cmax=5$, $p(error) = 6.9 \times 10^{-6}$; when $pmin = 0.9$ and $cmax=4$, $p(error) = 1.6 \times 10^{-7}$.

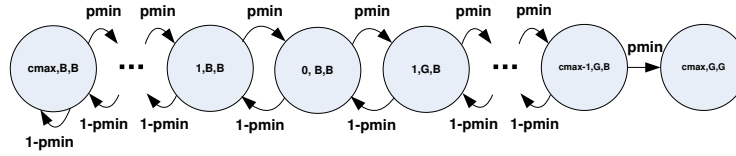


Fig. 5. State transition graph for convergence span

8.5 Proof of Theorem 5

By the definition of *convergence span*, the state transition graph of the stable system for convergence span is shown in Figure 5. The three elements in each global state are explained in the proof for theorem 4. From any state other than the start state (c_{max},B,B) and end state (c_{max},G,G) , the probability of going one step forward to the end state (c_{max},G,G) is p_{min} , while the probability of going one step backward towards the start state (c_{max},B,B) is $1-p_{min}$. So this is equivalent to say that from any state other than the end state (c_{max},G,G) and start state (c_{max},B,B) , the probability of going absolute one step forward is $p_{min}-(1-p_{min})$, which is $2 \times p_{min} - 1$, and the probability of staying in that state is $2 - 2 \times p_{min}$. For simplicity, we consider the approximate formula of the convergence span. That is, we make the probability of going one step forward from start state the same as that of other states, which is $2 \times p_{min} - 1$, instead of p_{min} . And accordingly the probability of staying in the start state is considered as $2 - 2 \times p_{min}$, instead of $1-p_{min}$. Since we need to go forward $2 \times c_{max}$ steps, with the probability of going one step forward as $2 \times p_{min} - 1$, the approximate average number of steps from start state to end state is $\frac{2 \times c_{max}}{2 \times p_{min} - 1}$.